

自底语法分析

2020年7月1日 21:13

1. 自底向上的语法分析：从分析树的底部(叶节点)向顶部(根节点)方向构造分析树
 - a. 可以看成是将输入串 w 归约为文法开始符号 S 的过程
 - b. 自顶向下的语法分析采用最左推导方式；自底向上的语法分析采用最左归约方式（反向构造最右推导）
 - c. 自底向上语法分析的通用框架：移入-归约分析(Shift-Reduce Parsing)
 - i. 在对输入串的一次从左到右扫描过程中，语法分析器将零个或多个输入符号移入到栈的顶端，直到可以对栈顶的一个文法符号串 β 进行归约为止
 - ii. 然后，它将 β 归约为某个产生式的左部
 - iii. 语法分析器不断地重复这个循环，直到它检测到一个语法错误，或者栈中包含了开始符号且输入缓冲区为空(当进入这样的格局时，语法分析器停止运行，并宣称成功完成了语法分析)为止
 - d. 移入-规约分析的动作
 - i. 移入：将下一个输入符号移到栈的顶端
 - ii. 归约：被归约的符号串的右端必然处于栈顶。语法分析器在栈中确定这个串的左端，并决定用哪个非终结符来替换这个串
 - iii. 接收：宣布语法分析过程成功完成
 - iv. 报错：发现一个语法错误，并调用错误恢复子例程
 - e. 移入-规约分析的特点
 - i. 若栈从左往右为从底到顶，输入从左到右为从顶到底，则 栈+输入 永远=规范句型
 - ii. 句柄：每次规约的符号串，应当是句型的最左直接短语，否则容易出错
2. LR分析法：从左扫描输入，反向构造出最右推导序列
 - a. LR文法(Knuth, 1963)是最大的可构造出相应移入-归约语法分析器的文法类
 - b. LR(k)分析：需要向前查看 k 个输入符号的LR分析
 - i. $k=0$ 和 $k=1$ 这两种情况具有实践意义，当省略(k)时，默认 $k=1$
 - c. 基本原理：逐步形成句柄，用移进 s 、待约、规约 r 等“状态”表示句柄识别进展程度
 - d. LR分析器（自动机）中除了符号栈还需要状态栈，对应的，分析表也需要（针对终结符的）动作表和（针对非终结符的）状态转移表，分析表中每行是当前状态，每列是处理中的符号
 - i. 初始化：状态栈： s_0 。符号栈： $\$$ 。输入区： $a_1a_2\dots a_n\$$
 - ii. 平常态： $s_0s_1\dots s_m$ 。 $\$X_1\dots X_m$ 。 $a_ia_{i+1}\dots a_n\$$
 - iii. $\text{ACTION}[s_m, a_i] = \text{sx}$ 后： $s_0s_1\dots s_msx$ 。 $\$X_1\dots X_m a_i$ 。 $a_{i+1}\dots a_n\$$
 - iv. $\text{ACTION}[s_m, a_i] = rx$ 后(产生式 $xA \rightarrow X_{m-(k-1)}\dots X_m$)： $s_0s_1\dots s_{m-k}$ 。 $\$X_1\dots X_{m-k} A$ 。 $a_{i+1}\dots a_n\$$ （从ii.平常态转移过来）
 - 1) $\text{GOTO}[s_{m-k}, A] = y$ 后： $s_0\dots s_{m-k} y$ 。 $\$X_1\dots X_{m-k} A$ 。 $a_{i+1}\dots a_n\$$ （紧接着

规约后, 从iv.规约态转移过来)

v. 如果ACTION[sm, ai]=acc, 那么分析成功

vi. 如果ACTION[sm, ai]=err, 那么出现语法错误

输入: 串w和LR语法分析表, 该表描述了文法G的ACTION函数和GOTO函数

输出: 如果w在L(G)中, 则输出w的自底向上语法分析过程中的归约步骤; 否则给出一个错误指示

方法: 初始时, 语法分析器栈中的内容为初始状态s0, 输入缓冲区中的内容为w\$。然后, 语法分析器执行下面的程序:

令a为w\$的第一个符号;

while(1) { /* 永远重复*/

 令s是栈顶的状态;

 if (ACTION [s, a]= st) {

 将t压入栈中;

 令a为下一个输入符号;

 } else if (ACTION [s, a]=归约A→β) {

 从栈中弹出 | β | 个符号;

 将GOTO[t, A]压入栈中;

 输出产生式A→β;

 } else if (ACTION [s, a]=接受)

 break; /* 语法分析完成*/

 else

 调用错误恢复例程;

}

3. LR(0)分析

a. LR(0)项目/项目: 右部的某位置有源点的产生式称为该文法的一个LR(0)项目

i. 项目描述了句柄识别的状态

ii. 右部有k个符号的项目有k+1个项目

1) 例: $S \rightarrow bBB$ 的项目:

a) 移进项目: $S \rightarrow \cdot bBB$

b) 待约项目: $S \rightarrow b \cdot BB$ 和 $S \rightarrow bB \cdot B$

c) 规约项目: $S \rightarrow bBB \cdot$

2) 例: 产生式 $A \rightarrow \epsilon$ 只生成一个项目 $A \rightarrow \cdot$

b. 后继项目 (Successive Item): 同属于一个产生式的项目, 但圆点的位置只相差一个符号, 则称后者是前者的后继项目

i. $GOTO(I, X) = CLOSURE(\{A \rightarrow \alpha X \cdot \beta \mid A \rightarrow \alpha \cdot X \beta \in I\})$

ii. 例: $A \rightarrow \alpha \cdot X \beta$ 的后继项目是 $A \rightarrow \alpha X \cdot \beta$

c. 项目集闭包(Closure of Item Sets): 把等待相同非终结符或其对应首符的等价项目组成集合

i. $CLOSURE(I) = I \cup \{B \rightarrow \cdot \gamma \mid A \rightarrow \alpha \cdot B \beta \in CLOSURE(I), B \rightarrow \gamma \in P\}$

ii. 每个项目集闭包对应着自动机一个状态

d. 增广文法(Augmented Grammar): 如果G是一个以S为开始符号的文法, 则G的增广文法G'就是在G中加上新开始符号S'和产生式 $S' \rightarrow S$ 而得到的文法

i. 引入这个新的开始产生式的目的是使得文法开始符号仅出现在一个产生式的左边, 从而使得分析器只有一个接受状态

e. LR(0)自动机: 基于增广文法, 分析每个项目集闭包的后继项目, 生成分析表

i. LR(0)自动机的规范状态集/项集族 (Canonical LR(0) Collection):

$C = \{I_0\} \cup \{I \mid \text{存在 } J \in C, X \in V_N \cup V_T, I = GOTO(J, X)\}$

- ii. 每条项目集闭包之间的转移边对应分析表的一行action s或GOTO
- iii. 对应增广文法的新产生式的项目集闭包所对应行只有一个表项，是\$列的action (acc)
- iv. 其他每个没有后继项目的项目集闭包在分析表中没有GOTO，但每个action都是r该项目集对应的产生式
- v. 给定文法 $G = (VN, VT, P, S)$ ，对应LR(0)自动机 $M = (C, V_NUV_T, GOTO, I0, F)$ 。其中状态集 $C = \{I0\} \cup \{I \mid \text{存在 } J \in C, X \in V_NUV_T, I = GOTO(J, X)\}$ ；初始状态 $I0 = CLOSURE(\{S' \rightarrow \cdot S\})$ ；终止状态 $F = \{CLOSURE(\{S' \rightarrow S \cdot\})\}$

SetOfItems CLOSURE(I) { //计算给定项目集I的闭包

```

J=I;
repeat
    for ( J中的每个项  $A \rightarrow \alpha \cdot B\beta$  )
        for ( G的每个产生式  $B \rightarrow \gamma$  )
            if (项  $B \rightarrow \cdot \gamma$  不在J中)
                将  $B \rightarrow \cdot \gamma$  加入J中;
until 在某一轮中没有新的项被加入到J中;
return J;

```

}

SetOfItems GOTO(I, X) { //返回项目集I对应于文法符号X的后继项目集闭包

```

将J初始化为空集;
for ( I中的每个项  $A \rightarrow \alpha \cdot X\beta$  )
    将项  $A \rightarrow \alpha X \cdot \beta$  加入到集合J中;
return CLOSURE (J);

```

}

void items(G') { //求项目集闭包

```

C = { CLOSURE ( {[S' → · S] } ) };
repeat
    for(C中的每个项集I)
        for(每个文法符号 X )
            if( GOTO( I, X)非空且不在C中)
                将GOTO( I, X)加入C中;
until在某一轮中没有新的项集被加入到C中;

```

}

//LR(0)分析表构造算法

构造 G' 的规范LR(0)项集族 $C = \{I0, I1, \dots, In\}$

令Ii对应状态i。状态i的语法分析动作按照下面的方法决定：

if $A \rightarrow \alpha \cdot a\beta \in Ii$ and $GOTO(Ii, a) = Ij$ then $ACTION[i, a] = sj$

if $A \rightarrow \alpha \cdot B\beta \in Ii$ and $GOTO(Ii, B) = Ij$ then $GOTO[i, B] = j$

if $A \rightarrow \alpha \cdot \in Ii$ and $A \neq S'$ then for $a \in V_T \cup \{\$ \}$ do $ACTION[i, a] = rj$ (j是产生式 $A \rightarrow \alpha$ 的编号)

if $S' \rightarrow S \cdot \in Ii$ then $ACTION[i, \$] = acc$

没有定义的所有条目都设置为“error”

//若按上述算法判断action时出现移进/规约冲突或规约/规约冲突（即一个元素多种动作），说明该CFG不是LR(0)文法

4. SLR分析（Simple的LR冲突化解方案）

a. 解决冲突的思想

- i. 设移进项目形如 $Ai \rightarrow \alpha i \cdot a_i \beta_i$ ；规约项目形如 $Bi \rightarrow \gamma i \cdot$ 。若 $FOLLOW(Bi)$ 两两不相交，且都不包含任一 a_i ，则按以下原则解决：
- ii. 输入符号a是 a_i 中的一个，则移进a；
- iii. 输入符号a属于 $FOLLOW(Bi)$ ，则规约成 $Bi \rightarrow \gamma i$ ；

iv. 否则报错

b. 构造算法:

构造G'的规范LR(0)项集族 $C = \{ I_0, I_1, \dots, I_n \}$

根据I构造得到状态i。状态i的语法分析动作按照下面的方法决定:

if $A \rightarrow \alpha \cdot a \beta \in I_i$ and $GOTO(I_i, a) = I_j$ then $ACTION[i, a] = sj$

if $A \rightarrow \alpha \cdot B \beta \in I_i$ and $GOTO(I_i, B) = I_j$ then $GOTO[i, B] = j$

if $A \rightarrow \alpha \cdot \in I_i$ and $A \neq S'$ then for $a \in FOLLOW(A)$ do $ACTION[i, a] = rj$ (j是产生式 $A \rightarrow \alpha$ 的编号)

if $S' \rightarrow S \cdot \in I_i$ then $ACTION[i, \$] = acc$;

没有定义的所有条目都设置为 "error"

//与传统LR区别是action[i,a]=rj的条件缩减了。若按上述算法判断action时出现移进/规约冲突或规约/规约冲突(即一个元素多种动作),说明该CFG不是SLR文法。(如FOLLOW集有运算符和\$时依旧不知道是移入还是规约好)

5. LR(1)分析

a. 解决冲突的思想

i. SLR中, 输入符号属于 $FOLLOW(B_i)$ 只是规约的必要条件, 而非充分条件

ii. LR(1)项: 形如 $[A \rightarrow \alpha \cdot \beta, a]$ 的项, 其中 $A \rightarrow \alpha \beta$ 是一个产生式, a是一个终结符(这里将\$也视为一个特殊的终结符)它表示在当前状态下, A后面必须紧跟的终结符, 称为该项的展望符(lookahead)

1) 当 β 不是 ϵ , 即**不在产生式最右时, 展望符并没有作用**

2) 当 β 是 ϵ , 即**在产生式最右时, 只有输入展望符a才能规约 $A \rightarrow \alpha$**

3) a的集合总是 $FOLLOW(A)$ 的子集, 而且一般是真子集

iii. 等价LR(1)项: 形如 $A \rightarrow \alpha \cdot B \beta, a$ 的项, 有等价项 $B \rightarrow \gamma \cdot b$

1) 其中b与产生式 $B \rightarrow \gamma$ 无关, 只与 β 和a有关, 即 $b \in FIRST(\beta a)$

2) 当 $\beta = \epsilon$ 时, $b = a$ 称作继承的后继符, 否则叫自生的后继符

项目集闭包 $CLOSURE(I) = I \cup \{ [B \rightarrow \gamma \cdot b] \mid [A \rightarrow \alpha \cdot B \beta, a] \in CLOSURE(I), B \rightarrow \gamma \in P, b \in FIRST(\beta a) \}$

SetOfItems $CLOSURE(I) \{$

repeat

for (I中的每个项 $[A \rightarrow \alpha \cdot B \beta, a]$

for (G'的每个产生式 $B \rightarrow \gamma$)

for ($FIRST(\beta a)$ 中的每个符号b)

将 $[B \rightarrow \gamma \cdot b]$ 加入到集合I中;

until不能向I中加入更多的项;

until I ;

}

后续集闭包 $GOTO(I, X) = CLOSURE(\{ [A \rightarrow \alpha X \cdot \beta, a] \mid [A \rightarrow \alpha \cdot X \beta, a] \in I \})$

SetOfItems $GOTO(I, X) \{$

将J初始化为空集;

for (I中的每个项 $[A \rightarrow \alpha \cdot X \beta, a]$)

将项 $[A \rightarrow \alpha X \cdot \beta, a]$ 加入到集合J中;

return $CLOSURE(J)$;

}

构造文法G'的LR(1)项集族

void items(G') {

将C初始化为 $\{ CLOSURE(\{ [S' \rightarrow \cdot S, \$] \}) \}$;

repeat

for (C中的每个项集I)

for (每个文法符号X)

if ($GOTO(I, X)$ 非空且不在C中)

将 $GOTO(I, X)$ 加入C中;

until 不再有新的项集加入到C中;

}

LR分析表构造算法

构造G'的规范LR(1)项集族 $C = \{I_0, I_1, \dots, I_n\}$

根据 I_i 构造得到状态 i 。状态 i 的语法分析动作按照下面的方法决定：

if $[A \rightarrow \alpha \cdot a\beta, b] \in I_i$ and $GOTO(I_i, a) = I_j$ then $ACTION[i, a] = sj$

if $[A \rightarrow \alpha \cdot B\beta, b] \in I_i$ and $GOTO(I_i, B) = I_j$ then $GOTO[i, B] = j$

if $[A \rightarrow \alpha \cdot, a] \in I_i$ and $A \neq S'$ then $ACTION[i, a] = rj$ (j 是产生式 $A \rightarrow \alpha$ 的编号)

if $[S' \rightarrow S \cdot, \$] \in I_i$ then $ACTION[i, \$] = acc$;

没有定义的所有条目都设置为“error”

若按上述算法判断action时出现移进/规约冲突或规约/规约冲突（即一个元素多种动作），说明该CFG不是SLR文法

6. LALR(1)分析 (LookAhead-LR)

a. 节省状态数的思路

- 同心项目集：除了展望符外都相同的项目集
- 将同心项目集合并（指让展望符变成展望符集合，常用/分割），若分析表没有冲突，称该文法为LALR(1)文法
- 合并后的展望符集合仍为FOLLOW集的子集

b. 冲突

- 合并后可能产生规约-规约冲突（指合并前，不同状态里同样的符号栈遇到同样的输入符可能使用不同产生式规约）
- 合并后不会产生移进-规约冲突（因为展望符只和规约有关）
- 合并后可能会使错误的发现被推迟（因为LALR可能会被合并后的展望符指引到多余的规约动作）

c. 特点

- 形式上与LR(1)相同
- 规模上和LR(0)/SLR相当
- 分析能力： $SLR < LALR(1) < LR(1)$ （有错，但错误不多）

7. 二义性文法的LR分析

- 二义性文法=有冲突的文法=不是LR文法
- 因为消除二义性要引入新的产生式，所以二义性文法一般更简短、自然
- 利用运算符优先级和结合性可解决一些冲突（如加号乘号，如if else）
- 必须在严格控制下才能保守使用二义性文法，因为稍有不慎就会让语法分析器所识别的语言出现偏差

8. LR分析的错误处理

- 可能的出错时机：查表发现空项
- 恐慌模式错误恢复：从栈顶向下扫描，直到发现某个状态 s_i ，它有一个对应于某个非终结符 A 的GOTO目标，可以认为从这个 A 推导出的串中包含错误
 - 然后丢弃0个或多个输入符号，直到发现一个可能合法地跟在 A 之后的符号 a 为止
 - 之后将 $s_{i+1} = GOTO(s_i, A)$ 压入栈中，继续进行正常的语法分析
- 短语层次错误恢复：检查LR分析表中的每一个报错条目，并根据语言的使用方法来决定程序员所犯的何种错误最有可能引起这个语法错误，然后构造出适当的恢复过程（将分析表的空白项填成对应的语法错误信息）

